

Advanced Machine Learning

Structured Models: Hidden Markov Models versus Conditional Random Fields



Eric Xing

Lecture 9, August 12, 2009

Reading:

Eric Xing

© Eric Xing @ CMU, 2006-2009

1



Dynamic clustering



Eric Xing

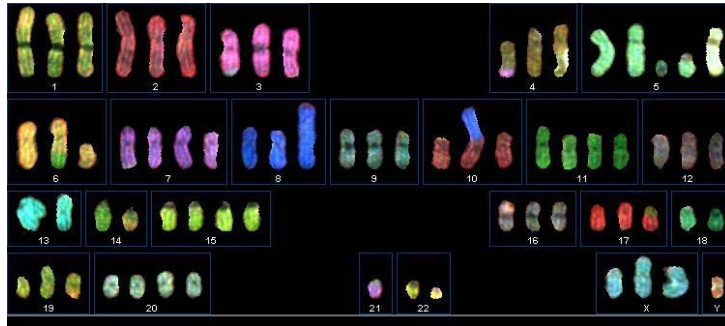
© Eric Xing @ CMU, 2006-2009

2

Clustering Tumor Cell States



Chromosomes of tumor cell:

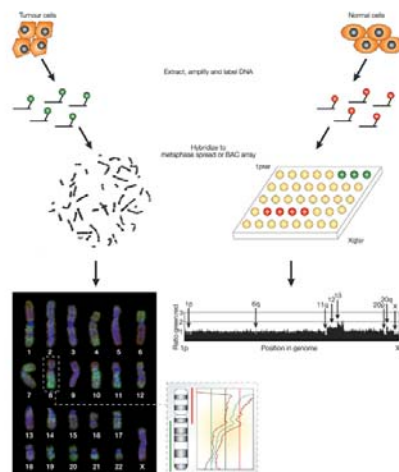


Eric Xing

© Eric Xing @ CMU, 2006-2009

3

Array CGH (comparative genomic hybridization)



- The basic assumption of a CGH experiment is that the ratio of the binding of test and control DNA is proportional to the ratio of the copy numbers of sequences in the two samples.
- But various kinds of noises make the true observations less easy to interpret ...

Eric Xing

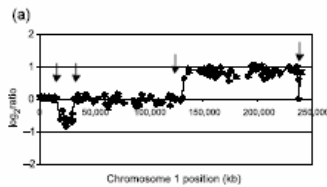
© Eric Xing @ CMU, 2006-2009

4

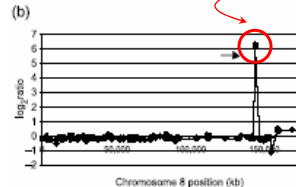
DNA Copy number aberration types in breast cancer



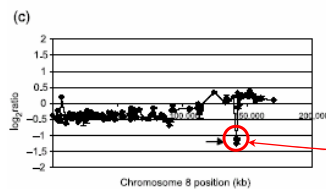
60-70 fold amplification of CMYC region



Copy number profile for chromosome 1 from 600 MPE cell line



Copy number profile for chromosome 8 from COLO320 cell line



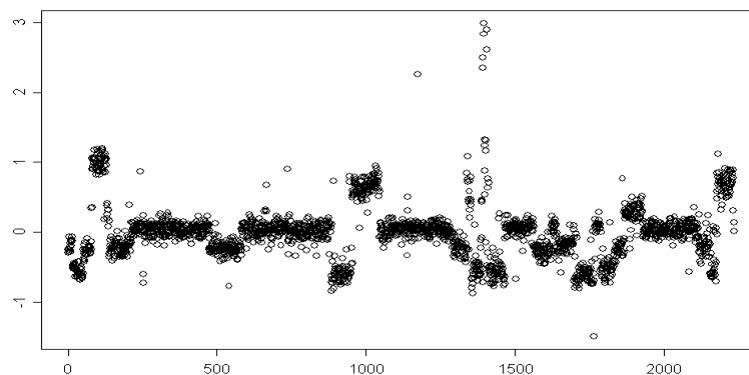
Copy number profile for chromosome 8 in MDA-MB-231 cell line

Eric Xing

© Eric Xing @ CMU, 2006-2009

5

A real CGH run



Eric Xing

© Eric Xing @ CMU, 2006-2009

6

Out problem: how to cluster sequential data?



Eric Xing

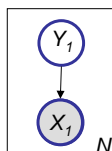
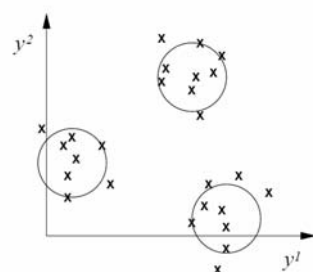
© Eric Xing @ CMU, 2006-2009

7

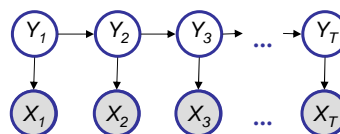
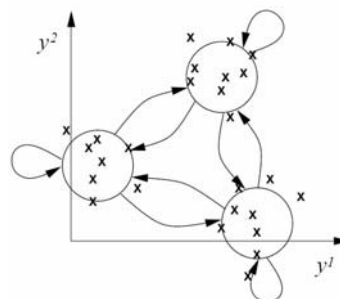
Hidden Markov Model: from static to dynamic mixture models



Static mixture



Dynamic mixture



Eric Xing

© Eric Xing @ CMU, 2006-2009

8

The Dishonest Casino

A casino has two dice:

- Fair die
 $P(1) = P(2) = P(3) = P(5) = P(6) = 1/6$
- Loaded die
 $P(1) = P(2) = P(3) = P(5) = 1/10$
 $P(6) = 1/2$

Casino player switches back-&-forth between fair and loaded die once every 20 turns

Game:

1. You bet \$1
2. You roll (always with a fair die)
3. Casino player rolls (maybe with fair die, maybe with loaded die)
4. Highest number wins \$2



Eric Xing

© Eric Xing @ CMU, 2006-2009

9

Definition (of HMM)

- Observation space

Alphabetic set: $C = \{c_1, c_2, \dots, c_K\}$
 Euclidean space: \mathbb{R}^d

- Index set of hidden states

$$I = \{1, 2, \dots, M\}$$

- Transition probabilities between any two states

$$p(y_t^j = 1 | y_{t-1}^i = 1) = a_{i,j},$$

$$\text{or } p(y_t | y_{t-1}^i = 1) \sim \text{Multinomial}(a_{i,1}, a_{i,2}, \dots, a_{i,M}), \forall i \in I.$$

- Start probabilities

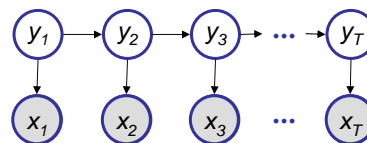
$$p(y_1) \sim \text{Multinomial}(\pi_1, \pi_2, \dots, \pi_M).$$

- Emission probabilities associated with each state

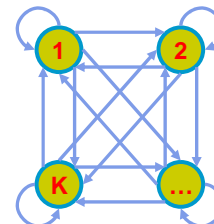
$$p(x_t | y_t^i = 1) \sim \text{Multinomial}(b_{i,1}, b_{i,2}, \dots, b_{i,K}), \forall i \in I.$$

or in general:

$$p(x_t | y_t^i = 1) \sim f(\cdot | \theta_i), \forall i \in I.$$



Graphical model



State automata

Eric Xing

© Eric Xing @ CMU, 2006-2009

10

Applications of HMMs



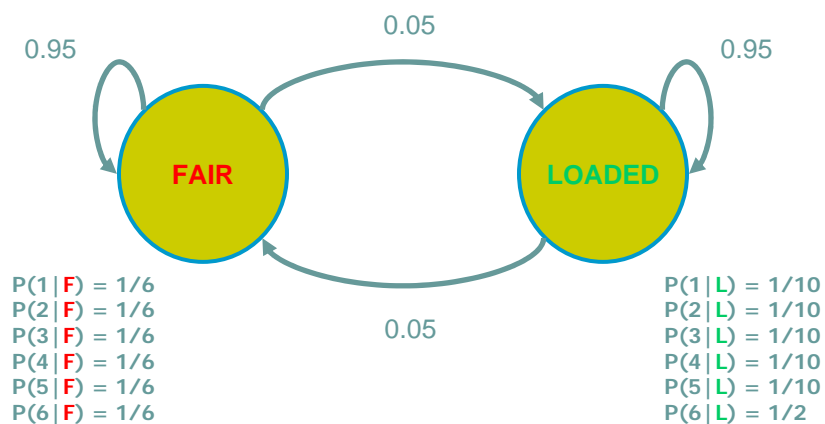
- **Some early applications of HMMs**
 - finance, but we never saw them
 - speech recognition
 - modelling ion channels
- In the mid-late 1980s HMMs entered genetics and molecular biology, and they are now firmly entrenched.
- **Some current applications of HMMs to biology**
 - mapping chromosomes
 - aligning biological sequences
 - predicting sequence structure
 - inferring evolutionary relationships
 - finding genes in DNA sequence

Eric Xing

© Eric Xing @ CMU, 2006-2009

11

The Dishonest Casino Model



Eric Xing

© Eric Xing @ CMU, 2006-2009

12

Puzzles Regarding the Dishonest Casino



GIVEN: A sequence of rolls by the casino player

1245526462146146136136661664661636616366163616515615115146123562344

QUESTION

- How likely is this sequence, given our model of how the casino works?
 - This is the **EVALUATION** problem in HMMs
- What portion of the sequence was generated with the fair die, and what portion with the loaded die?
 - This is the **DECODING** question in HMMs
- How “loaded” is the loaded die? How “fair” is the fair die? How often does the casino player change from fair to loaded, and back?
 - This is the **LEARNING** question in HMMs

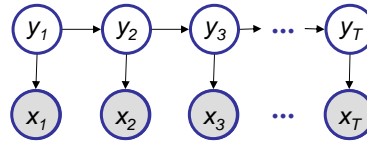
Joint Probability



1245526462146146136136661664661636616366163616515615115146123562344

Probability of a Parse

- Given a sequence $\mathbf{x} = x_1, \dots, x_T$ and a parse $\mathbf{y} = y_1, \dots, y_T$,
- To find how likely is the parse:
(given our HMM and the sequence)



$$\begin{aligned}
 p(\mathbf{x}, \mathbf{y}) &= p(x_1, \dots, x_T, y_1, \dots, y_T) && \text{(Joint probability)} \\
 &= p(y_1) p(x_1 | y_1) p(y_2 | y_1) p(x_2 | y_2) \dots p(y_T | y_{T-1}) p(x_T | y_T) \\
 &= p(y_1) P(y_2 | y_1) \dots p(y_T | y_{T-1}) \times p(x_1 | y_1) p(x_2 | y_2) \dots p(x_T | y_T)
 \end{aligned}$$

- Marginal probability: $p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}) = \sum_{y_1} \sum_{y_2} \dots \sum_{y_N} \pi_{y_1} \prod_{t=2}^T a_{y_{t-1} y_t} \prod_{t=1}^T p(x_t | y_t)$
- Posterior probability: $p(\mathbf{y} | \mathbf{x}) = p(\mathbf{x}, \mathbf{y}) / p(\mathbf{x})$

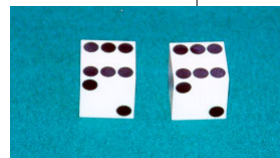
Eric Xing

© Eric Xing @ CMU, 2006-2009

15

Example: the Dishonest Casino

- Let the sequence of rolls be:
 - $\mathbf{x} = 1, 2, 1, 5, 6, 2, 1, 6, 2, 4$
- Then, what is the likelihood of
 - $\mathbf{y} = \text{Fair, Fair, Fair, Fair, Fair, Fair, Fair, Fair, Fair, Fair?}$
(say initial probs $a_{0\text{Fair}} = 1/2$, $a_{0\text{Loaded}} = 1/2$)



$$\frac{1}{2} \times P(1 | \text{Fair}) P(\text{Fair} | \text{Fair}) P(2 | \text{Fair}) P(\text{Fair} | \text{Fair}) \dots P(4 | \text{Fair}) =$$

$$\frac{1}{2} \times (1/6)^{10} \times (0.95)^9 = .00000000521158647211 = 5.21 \times 10^{-9}$$

Eric Xing

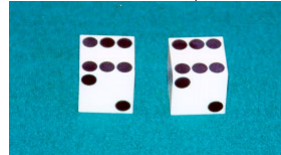
© Eric Xing @ CMU, 2006-2009

16

Example: the Dishonest Casino



- So, the likelihood the die is fair in all this run is just 5.21×10^{-9}



- OK, but what is the likelihood of
 - π = Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded?

$$\frac{1}{2} \times P(1 \mid \text{Loaded}) P(\text{Loaded} \mid \text{Loaded}) \dots P(4 \mid \text{Loaded}) =$$

$$\frac{1}{2} \times (1/10)^8 \times (1/2)^2 (0.95)^9 = .00000000078781176215 = 0.79 \times 10^{-9}$$

- Therefore, it is after all 6.59 times more likely that the die is fair all the way, than that it is loaded all the way

Eric Xing

© Eric Xing @ CMU, 2006-2009

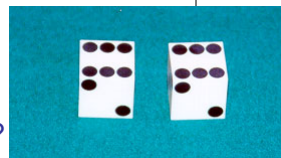
17

Example: the Dishonest Casino



- Let the sequence of rolls be:

- $x = 1, 6, 6, 5, 6, 2, 6, 6, 3, 6$



- Now, what is the likelihood $\pi = F, F, \dots, F$?

- $\frac{1}{2} \times (1/6)^{10} \times (0.95)^9 = 0.5 \times 10^{-9}$, same as before

- What is the likelihood $y = L, L, \dots, L$?

$$\frac{1}{2} \times (1/10)^4 \times (1/2)^6 (0.95)^9 = .00000049238235134735 = 5 \times 10^{-7}$$

- So, it is 100 times more likely the die is loaded

Eric Xing

© Eric Xing @ CMU, 2006-2009

18

Three Main Questions on HMMs



1. Evaluation

GIVEN an HMM \mathcal{M} , and a sequence \mathbf{x} ,
 FIND Prob ($\mathbf{x} | \mathcal{M}$)
 ALGO. **Forward**

2. Decoding

GIVEN an HMM \mathcal{M} , and a sequence \mathbf{x} ,
 FIND the sequence \mathbf{y} of states that maximizes, e.g., $P(\mathbf{y} | \mathbf{x}, \mathcal{M})$,
 or the most probable subsequence of states
 ALGO. **Viterbi, Forward-backward**

3. Learning

GIVEN an HMM \mathcal{M} , with unspecified transition/emission probs.,
 and a sequence \mathbf{x} ,
 FIND parameters $\theta = (\pi_i, a_{ij}, \eta_{ik})$ that maximize $P(\mathbf{x} | \theta)$
 ALGO. **Baum-Welch (EM)**

The Forward Algorithm



- We want to calculate $P(\mathbf{x})$, the likelihood of \mathbf{x} , given the HMM
 - Sum over all possible ways of generating \mathbf{x} :

$$p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}) = \sum_{y_1} \sum_{y_2} \cdots \sum_{y_N} \pi_{y_1} \prod_{t=2}^T a_{y_{t-1}, y_t} \prod_{t=1}^T p(x_t | y_t)$$

- To avoid summing over an exponential number of paths \mathbf{y} , define

$$\alpha(y_t^k = 1) = \alpha_t^k \stackrel{\text{def}}{=} P(x_1, \dots, x_t, y_t^k = 1) \quad (\text{the forward probability})$$

- The recursion:

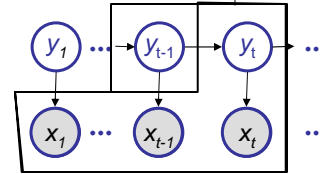
$$\alpha_t^k = p(x_t | y_t^k = 1) \sum_i \alpha_{t-1}^i a_{i,k}$$

$$P(\mathbf{x}) = \sum_k \alpha_T^k$$

The Forward Algorithm – derivation



- Compute the forward probability:



$$\alpha_t^k = P(x_1, \dots, x_{t-1}, x_t, y_t^k = 1)$$

$$\begin{aligned} &= \sum_{y_{t-1}} P(x_1, \dots, x_{t-1}, y_{t-1}) P(y_t^k = 1 | y_{t-1}, x_1, \dots, x_{t-1}) P(x_t | y_t^k = 1, x_1, \dots, x_{t-1}, y_{t-1}) \\ &= \sum_{y_{t-1}} P(x_1, \dots, x_{t-1}, y_{t-1}) P(y_t^k = 1 | y_{t-1}) P(x_t | y_t^k = 1) \\ &= P(x_t | y_t^k = 1) \sum_i P(x_1, \dots, x_{t-1}, y_{t-1}^i = 1) P(y_t^k = 1 | y_{t-1}^i = 1) \\ &= P(x_t | y_t^k = 1) \sum_i \alpha_{t-1}^i a_{i,k} \end{aligned}$$

Chain rule: $P(A, B, C) = P(A)P(B | A)P(C | A, B)$

Eric Xing

© Eric Xing @ CMU, 2006-2009

21

The Forward Algorithm



- We can compute α_t^k for all k, t , using dynamic programming!

Initialization:

$$\alpha_1^k = P(x_1 | y_1^k = 1) \pi_k$$

$$\begin{aligned} \alpha_1^k &= P(x_1, y_1^k = 1) \\ &= P(x_1 | y_1^k = 1) P(y_1^k = 1) \\ &= P(x_1 | y_1^k = 1) \pi_k \end{aligned}$$

Iteration:

$$\alpha_t^k = P(x_t | y_t^k = 1) \sum_i \alpha_{t-1}^i a_{i,k}$$

Termination:

$$P(\mathbf{x}) = \sum_k \alpha_T^k$$

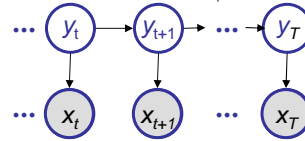
Eric Xing

© Eric Xing @ CMU, 2006-2009

22

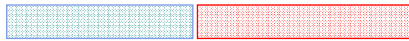
The Backward Algorithm

- We want to compute $P(y_t^k = 1 | \mathbf{x})$,
the posterior probability distribution on the t^{th} position, given \mathbf{x}



- We start by computing

$$\begin{aligned} P(y_t^k = 1, \mathbf{x}) &= P(x_1, \dots, x_t, y_t^k = 1, x_{t+1}, \dots, x_T) \\ &= P(x_1, \dots, x_t, y_t^k = 1) P(x_{t+1}, \dots, x_T | x_1, \dots, x_t, y_t^k = 1) \\ &= P(x_1 \dots x_t, y_t^k = 1) P(x_{t+1} \dots x_T | y_t^k = 1) \end{aligned}$$



Forward, α_t^k

Backward, $\beta_t^k = P(x_{t+1}, \dots, x_T | y_t^k = 1)$

- The recursion:

$$\beta_t^k = \sum_i a_{k,i} p(x_{t+1} | y_{t+1}^i = 1) \beta_{t+1}^i$$

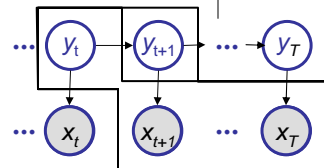
Eric Xing

© Eric Xing @ CMU, 2006-2009

23

The Backward Algorithm – derivation

- Define the backward probability:



$$\begin{aligned} \beta_t^k &= P(x_{t+1}, \dots, x_T | y_t^k = 1) \\ &= \sum_{y_{t+1}^i} P(x_{t+1}, \dots, x_T, y_{t+1}^i | y_t^k = 1) \\ &= \sum_i P(y_{t+1}^i = 1 | y_t^k = 1) p(x_{t+1} | y_{t+1}^i = 1, y_t^k = 1) P(x_{t+2}, \dots, x_T | x_{t+1}, y_{t+1}^i = 1, y_t^k = 1) \\ &= \sum_i P(y_{t+1}^i = 1 | y_t^k = 1) p(x_{t+1} | y_{t+1}^i = 1) P(x_{t+2}, \dots, x_T | y_{t+1}^i = 1) \\ &= \sum_i a_{k,i} p(x_{t+1} | y_{t+1}^i = 1) \beta_{t+1}^i \end{aligned}$$

Chain rule: $P(A, B, C | \alpha) = P(A | \alpha) P(B | A, \alpha) P(C | A, B, \alpha)$

Eric Xing

© Eric Xing @ CMU, 2006-2009

24

The Backward Algorithm



- We can compute β_t^k for all k, t , using dynamic programming!

Initialization:

$$\beta_T^k = 1, \forall k$$

Iteration:

$$\beta_t^k = \sum_i a_{k,i} P(x_{t+1} | y_{t+1}^i = 1) \beta_{t+1}^i$$

Termination:

$$P(x) = \sum_k \alpha_1^k \beta_1^k$$

Eric Xing

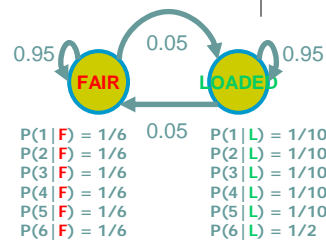
© Eric Xing @ CMU, 2006-2009

25

Example:



$x = 1, 2, 1, 5, 6, 2, 1, 6, 2, 4$



$$\alpha_t^k = P(x_t | y_t^k = 1) \sum_i \alpha_{t-1}^i a_{i,k}$$

$$\beta_t^k = \sum_i a_{k,i} P(x_{t+1} | y_{t+1}^i = 1) \beta_{t+1}^i$$

Eric Xing

© Eric Xing @ CMU, 2006-2009

26

Posterior decoding

- We can now calculate

$$P(y_t^k = 1 | \mathbf{x}) = \frac{P(y_t^k = 1, \mathbf{x})}{P(\mathbf{x})} = \frac{\alpha_t^k \beta_t^k}{P(\mathbf{x})}$$

- Then, we can ask

- What is the most likely state at position t of sequence \mathbf{x} :

$$k_t^* = \arg \max_k P(y_t^k = 1 | \mathbf{x})$$

- Note that this is an MPA of a **single** hidden state, what if we want to a MPA of a whole hidden state sequence?

- Posterior Decoding: $\{y_t^{k_t^*} = 1 : t = 1 \dots T\}$

- This is different from MPA of a **whole** sequence states

- This can be understood as **bit error rate** vs. **word error rate**

Example:
MPA of X ?
MPA of (X, Y) ?

		of hidden
x	y	$P(x, y)$
0	0	0.35
0	1	0.05
1	0	0.3
1	1	0.3

Eric Xing

© Eric Xing @ CMU, 2006-2009

27

Viterbi decoding

- GIVEN $\mathbf{x} = x_1, \dots, x_T$, we want to find $\mathbf{y} = y_1, \dots, y_T$, such that $P(\mathbf{y} | \mathbf{x})$ is maximized:

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} P(\mathbf{y} | \mathbf{x}) = \operatorname{argmax}_{\pi} P(\mathbf{y}, \mathbf{x})$$

- Let

$$V_t^k = \max_{\{y_1, \dots, y_{t-1}\}} P(x_1, \dots, x_{t-1}, y_1, \dots, y_{t-1}, x_t, y_t^k = 1)$$

= Probability of most likely **sequence of states** ending at state $y_t = k$

- The recursion:

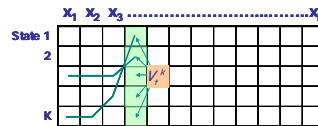
$$V_t^k = p(x_t | y_t^k = 1) \max_i a_{i,k} V_{t-1}^i$$

- Underflows are a significant problem

$$p(x_1, \dots, x_t, y_1, \dots, y_t) = \pi_{y_1} a_{y_1, y_2} \dots a_{y_{t-1}, y_t} b_{y_1, x_1} \dots b_{y_t, x_t}$$

- These numbers become extremely small – underflow

- Solution: Take the logs of all values: $V_t^k = \log p(x_t | y_t^k = 1) + \max_i (\log(a_{i,k}) + V_{t-1}^i)$



Eric Xing

© Eric Xing @ CMU, 2006-2009

28

Computational Complexity and implementation details



- What is the running time, and space required, for Forward, and Backward?

$$\alpha_t^k = p(x_t | y_t^k = 1) \sum_i \alpha_{t-1}^i a_{i,k}$$

$$\beta_t^k = \sum_i a_{k,i} p(x_{t+1} | y_{t+1}^i = 1) \beta_{t+1}^i$$

$$V_t^k = p(x_t | y_t^k = 1) \max_i a_{i,k} V_{t-1}^i$$

Time: $O(K^2M)$; Space: $O(KM)$.

- Useful implementation technique to avoid underflows
 - Viterbi: sum of logs
 - Forward/Backward: rescaling at each position by multiplying by a constant

Eric Xing

© Eric Xing @ CMU, 2006-2009

29

Learning HMM: two scenarios



- Supervised learning:** estimation when the “right answer” is known
 - Examples:**
 - GIVEN:** a genomic region $x = x_1 \dots x_{1,000,000}$ where we have good (experimental) annotations of the CpG islands
 - GIVEN:** the casino player allows us to observe him one evening, as he changes dice and produces 10,000 rolls
- Unsupervised learning:** estimation when the “right answer” is unknown
 - Examples:**
 - GIVEN:** the porcupine genome; we don't know how frequent are the CpG islands there, neither do we know their composition
 - GIVEN:** 10,000 rolls of the casino player, but we don't see when he changes dice
- QUESTION:** Update the parameters θ of the model to maximize $P(x|\theta)$ --- Maximal likelihood (ML) estimation

Eric Xing

© Eric Xing @ CMU, 2006-2009

30

Supervised ML estimation



- Given $\mathbf{x} = x_1 \dots x_N$ for which the true state path $\mathbf{y} = y_1 \dots y_N$ is known,

- Define:

$$\begin{aligned} A_{ij} &= \# \text{ times state transition } i \rightarrow j \text{ occurs in } \mathbf{y} \\ B_{ik} &= \# \text{ times state } i \text{ in } \mathbf{y} \text{ emits } k \text{ in } \mathbf{x} \end{aligned}$$

- We can show that the **maximum likelihood** parameters θ are:

$$a_{ij}^{ML} = \frac{\#(i \rightarrow j)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=2}^T y_{n,t-1}^i y_{n,t}^j}{\sum_n \sum_{t=2}^T y_{n,t-1}^i} = \frac{A_{ij}}{\sum_j A_{ij}}$$

$$b_{ik}^{ML} = \frac{\#(i \rightarrow k)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=1}^T y_{n,t}^i x_{n,t}^k}{\sum_n \sum_{t=1}^T y_{n,t}^i} = \frac{B_{ik}}{\sum_k B_{ik}}$$

(Homework!)

- What if \mathbf{y} is continuous? We can treat $\{(x_{n,t}, y_{n,t}) : t=1:T, n=1:N\}$ as $N \times T$ observations of, e.g., a Gaussian, and apply learning rules for Gaussian ...

(Homework!)

Eric Xing

© Eric Xing @ CMU, 2006-2009

31

Pseudocounts



- Solution for small training sets:

- Add pseudocounts

$$\begin{aligned} A_{ij} &= \# \text{ times state transition } i \rightarrow j \text{ occurs in } \mathbf{y} + R_{ij} \\ B_{ik} &= \# \text{ times state } i \text{ in } \mathbf{y} \text{ emits } k \text{ in } \mathbf{x} + S_{ik} \end{aligned}$$

- R_{ij}, S_{ij} are pseudocounts representing our prior belief

- Total pseudocounts: $R_i = \sum_j R_{ij}, S_i = \sum_k S_{ik}$,

- "strength" of prior belief,
- total number of imaginary instances in the prior

- Larger total pseudocounts \Rightarrow **strong prior belief**

- Small total pseudocounts: just to avoid 0 probabilities --- **smoothing**

Eric Xing

© Eric Xing @ CMU, 2006-2009

32

Unsupervised ML estimation

- Given $\mathbf{x} = x_1 \dots x_N$ for which the true state path $y = y_1 \dots y_N$ is unknown,

- EXPECTATION MAXIMIZATION**

- Starting with our best guess of a model \mathcal{M} , parameters θ .
- Estimate A_{ij}, B_{ik} in the training data
 - How? $A_{ij} = \sum_{n,t} \langle y_{n,t-1}^i y_{n,t}^j \rangle$ $B_{ik} = \sum_{n,t} \langle y_{n,t}^i x_{n,t}^k \rangle$, How? (homework)
- Update θ according to A_{ij}, B_{ik}
 - Now a "supervised learning" problem
- Repeat 1 & 2, until convergence

This is called the Baum-Welch Algorithm

We can get to a provably more (or equally) likely parameter set θ each iteration

Eric Xing

© Eric Xing @ CMU, 2006-2009

33

The Baum Welch algorithm

- The complete log likelihood

$$\ell_c(\theta; \mathbf{x}, \mathbf{y}) = \log p(\mathbf{x}, \mathbf{y}) = \log \prod_n \left(p(y_{n,1}) \prod_{t=2}^T p(y_{n,t} | y_{n,t-1}) \prod_{t=1}^T p(x_{n,t} | y_{n,t}) \right)$$

- The expected complete log likelihood

$$\langle \ell_c(\theta; \mathbf{x}, \mathbf{y}) \rangle = \sum_n \left(\langle y_{n,1}^i \rangle_{p(y_{n,1} | \mathbf{x}_n)} \log \pi_i \right) + \sum_n \sum_{t=2}^T \left(\langle y_{n,t-1}^i y_{n,t}^j \rangle_{p(y_{n,t-1}, y_{n,t} | \mathbf{x}_n)} \log a_{i,j} \right) + \sum_n \sum_{t=1}^T \left(\langle x_{n,t}^k y_{n,t}^i \rangle_{p(y_{n,t} | \mathbf{x}_n)} \log b_{i,k} \right)$$

- EM

- The E step

$$\gamma_{n,t}^i = \langle y_{n,t}^i \rangle = p(y_{n,t}^i = 1 | \mathbf{x}_n)$$

$$\xi_{n,t}^{i,j} = \langle y_{n,t-1}^i y_{n,t}^j \rangle = p(y_{n,t-1}^i = 1, y_{n,t}^j = 1 | \mathbf{x}_n)$$

- The M step ("symbolically" identical to MLE)

$$\pi_i^{ML} = \frac{\sum_n \gamma_{n,1}^i}{N}$$

$$a_{ij}^{ML} = \frac{\sum_n \sum_{t=2}^T \xi_{n,t}^{i,j}}{\sum_n \sum_{t=1}^{T-1} \gamma_{n,t}^i}$$

$$b_{ik}^{ML} = \frac{\sum_n \sum_{t=1}^T \gamma_{n,t}^i x_{n,t}^k}{\sum_n \sum_{t=1}^T \gamma_{n,t}^i}$$

Eric Xing

© Eric Xing @ CMU, 2006-2009

34

Summary



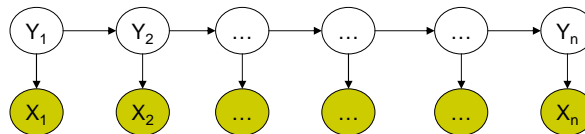
- Modeling hidden transitional trajectories (in discrete state space, such as cluster label, DNA copy number, dice-choice, etc.) underlying observed sequence data (discrete, such as dice outcomes; or continuous, such as CGH signals)
- Useful for parsing, segmenting sequential data
- Important HMM computations:
 - The joint likelihood of a parse and data can be written as a product to local terms (i.e., initial prob, transition prob, emission prob.)
 - Computing marginal likelihood of the observed sequence: forward algorithm
 - Predicting a single hidden state: forward-backward
 - Predicting an entire sequence of hidden states: viterbi
 - Learning HMM parameters: an EM algorithm known as Baum-Welch

Eric Xing

© Eric Xing @ CMU, 2006-2009

35

Shortcomings of Hidden Markov Model



- HMM models capture dependences between each state and **only** its corresponding observation
 - NLP example: In a sentence segmentation task, each segmental state may depend not just on a single word (and the adjacent segmental stages), but also on the (non-local) features of the whole line such as line length, indentation, amount of white space, etc.
- Mismatch between learning objective function and prediction objective function
 - HMM learns a joint distribution of states and observations $P(Y, X)$, but in a prediction task, we need the conditional probability $P(Y|X)$

Eric Xing

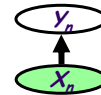
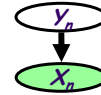
© Eric Xing @ CMU, 2006-2009

36

Recall Generative vs. Discriminative Classifiers



- Goal: Wish to learn $f: X \rightarrow Y$, e.g., $P(Y|X)$
- Generative classifiers (e.g., Naïve Bayes):
 - Assume some functional form for $P(X|Y)$, $P(Y)$
This is a '**generative**' model of the data!
 - Estimate parameters of $P(X|Y)$, $P(Y)$ directly from training data
 - Use Bayes rule to calculate $P(Y|X=x)$
- Discriminative classifiers (e.g., logistic regression)
 - Directly assume some functional form for $P(Y|X)$
This is a '**discriminative**' model of the data!
 - Estimate parameters of $P(Y|X)$ directly from training data

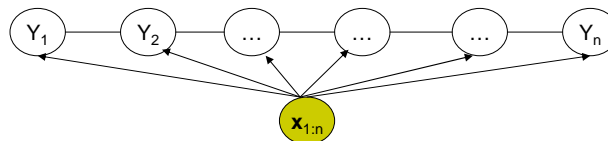


Eric Xing

© Eric Xing @ CMU, 2006-2009

37

Conditional Random Fields



$$P(\mathbf{y}_{1:n} | \mathbf{x}_{1:n}) = \frac{1}{Z(\mathbf{x}_{1:n})} \prod_{i=1}^n \phi(y_i, y_{i-1}, \mathbf{x}_{1:n}) = \frac{1}{Z(\mathbf{x}_{1:n}, \mathbf{w})} \prod_{i=1}^n \exp(\mathbf{w}^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}_{1:n}))$$

- CRF is a partially directed model
 - Discriminative model
 - Usage of global normalizer $Z(\mathbf{x})$
 - Models the dependence between each state and the entire observation sequence

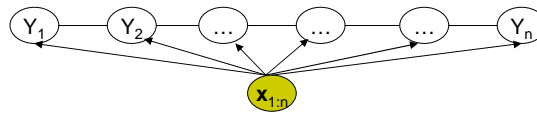
Eric Xing

© Eric Xing @ CMU, 2006-2009

38

Conditional Random Fields

- General parametric form:



$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x}, \lambda, \mu)} \exp\left(\sum_{i=1}^n \left(\sum_k \lambda_k f_k(y_i, y_{i-1}, \mathbf{x}) + \sum_l \mu_l g_l(y_i, \mathbf{x})\right)\right)$$

$$= \frac{1}{Z(\mathbf{x}, \lambda, \mu)} \exp\left(\sum_{i=1}^n (\lambda^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}) + \mu^T \mathbf{g}(y_i, \mathbf{x}))\right)$$

$$\text{where } Z(\mathbf{x}, \lambda, \mu) = \sum_{\mathbf{y}} \exp\left(\sum_{i=1}^n (\lambda^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}) + \mu^T \mathbf{g}(y_i, \mathbf{x}))\right)$$

Eric Xing

© Eric Xing @ CMU, 2006-2009

39

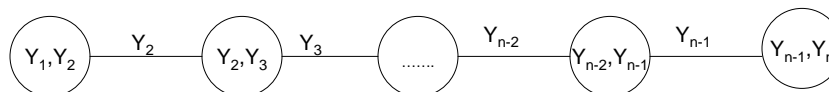
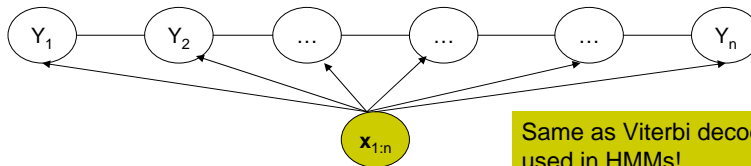
CRFs: Inference

- Given CRF parameters λ and μ , find the \mathbf{y}^* that maximizes $P(\mathbf{y}|\mathbf{x})$

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} \exp\left(\sum_{i=1}^n (\lambda^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}) + \mu^T \mathbf{g}(y_i, \mathbf{x}))\right)$$

- Can ignore $Z(\mathbf{x})$ because it is not a function of \mathbf{y}

- Run the max-product algorithm on the junction-tree of CRF:



Eric Xing

© Eric Xing @ CMU, 2006-2009

40

CRF learning

- Given $\{(\mathbf{x}_d, \mathbf{y}_d)\}_{d=1}^N$, find λ^*, μ^* such that

$$\begin{aligned}\lambda^*, \mu^* &= \arg \max_{\lambda, \mu} L(\lambda, \mu) = \arg \max_{\lambda, \mu} \prod_{d=1}^N P(\mathbf{y}_d | \mathbf{x}_d, \lambda, \mu) \\ &= \arg \max_{\lambda, \mu} \prod_{d=1}^N \frac{1}{Z(\mathbf{x}_d, \lambda, \mu)} \exp\left(\sum_{i=1}^n (\lambda^T \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d) + \mu^T \mathbf{g}(y_{d,i}, \mathbf{x}_d))\right) \\ &= \arg \max_{\lambda, \mu} \sum_{d=1}^N \left(\sum_{i=1}^n (\lambda^T \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d) + \mu^T \mathbf{g}(y_{d,i}, \mathbf{x}_d)) - \log Z(\mathbf{x}_d, \lambda, \mu) \right)\end{aligned}$$

Gradient of the log-partition function in an exponential family is the expectation of the sufficient statistics.

- Computing the gradient w.r.t λ :

$$\nabla_{\lambda} L(\lambda, \mu) = \sum_{d=1}^N \left(\sum_{i=1}^n \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d) - \sum_{\mathbf{y}} (P(\mathbf{y} | \mathbf{x}_d) \sum_{i=1}^n \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d)) \right)$$

Eric Xing

© Eric Xing @ CMU, 2006-2009

41

CRF learning

$$\nabla_{\lambda} L(\lambda, \mu) = \sum_{d=1}^N \left(\sum_{i=1}^n \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d) - \sum_{\mathbf{y}} (P(\mathbf{y} | \mathbf{x}_d) \sum_{i=1}^n \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d)) \right)$$

- Computing the model expectations:

- Requires exponentially large number of summations: Is it intractable?

$$\begin{aligned}\sum_{\mathbf{y}} (P(\mathbf{y} | \mathbf{x}_d) \sum_{i=1}^n \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d)) &= \sum_{i=1}^n \left(\sum_{\mathbf{y}} \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d) P(\mathbf{y} | \mathbf{x}_d) \right) \\ &= \sum_{i=1}^n \sum_{y_i, y_{i-1}} \mathbf{f}(y_i, y_{i-1}, \mathbf{x}_d) P(y_i, y_{i-1} | \mathbf{x}_d)\end{aligned}$$

Expectation of \mathbf{f} over the corresponding marginal probability of neighboring nodes!!

- Tractable!

- Can compute marginals using the sum-product algorithm on the chain

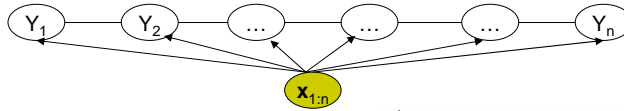
Eric Xing

© Eric Xing @ CMU, 2006-2009

42

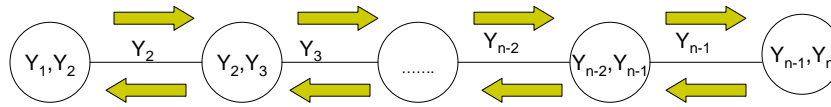
CRF learning

- Computing marginals using junction-tree calibration:



- Junction Tree Initialization:

$$\alpha^0(y_i, y_{i-1}) = \exp(\lambda^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}_d) + \mu^T \mathbf{g}(y_i, \mathbf{x}_d))$$



- After calibration:

$$P(y_i, y_{i-1} | \mathbf{x}_d) \propto \alpha(y_i, y_{i-1})$$

$$\Rightarrow P(y_i, y_{i-1} | \mathbf{x}_d) = \frac{\alpha(y_i, y_{i-1})}{\sum_{y_i, y_{i-1}} \alpha(y_i, y_{i-1})} = \alpha'(y_i, y_{i-1})$$

Also called forward-backward algorithm

Eric Xing

© Eric Xing @ CMU, 2006-2009

43

CRF learning

- Computing feature expectations using calibrated potentials:

$$\sum_{y_i, y_{i-1}} \mathbf{f}(y_i, y_{i-1}, \mathbf{x}_d) P(y_i, y_{i-1} | \mathbf{x}_d) = \sum_{y_i, y_{i-1}} \mathbf{f}(y_i, y_{i-1}, \mathbf{x}_d) \alpha'(y_i, y_{i-1})$$

- Now we know how to compute $\nabla_{\lambda} L(\lambda, \mu)$:

$$\begin{aligned} \nabla_{\lambda} L(\lambda, \mu) &= \sum_{d=1}^N \left(\sum_{i=1}^n \mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d) - \sum_{\mathbf{y}} (P(\mathbf{y} | \mathbf{x}_d) \sum_{i=1}^n \mathbf{f}(y_i, y_{i-1}, \mathbf{x}_d)) \right) \\ &= \sum_{d=1}^N \left(\sum_{i=1}^n (\mathbf{f}(y_{d,i}, y_{d,i-1}, \mathbf{x}_d) - \sum_{y_i, y_{i-1}} \alpha'(y_i, y_{i-1}) \mathbf{f}(y_i, y_{i-1}, \mathbf{x}_d)) \right) \end{aligned}$$

- Learning can now be done using gradient ascent:

$$\begin{aligned} \lambda^{(t+1)} &= \lambda^{(t)} + \eta \nabla_{\lambda} L(\lambda^{(t)}, \mu^{(t)}) \\ \mu^{(t+1)} &= \mu^{(t)} + \eta \nabla_{\mu} L(\lambda^{(t)}, \mu^{(t)}) \end{aligned}$$

Eric Xing

© Eric Xing @ CMU, 2006-2009

44

CRF learning

- In practice, we use a Gaussian Regularizer for the parameter vector to improve generalizability

$$\lambda^*, \mu^* = \arg \max_{\lambda, \mu} \sum_{d=1}^N \log P(\mathbf{y}_d | \mathbf{x}_d, \lambda, \mu) - \frac{1}{2\sigma^2} (\lambda^T \lambda + \mu^T \mu)$$

- In practice, gradient ascent has very slow convergence
 - Alternatives:
 - Conjugate Gradient method
 - Limited Memory Quasi-Newton Methods

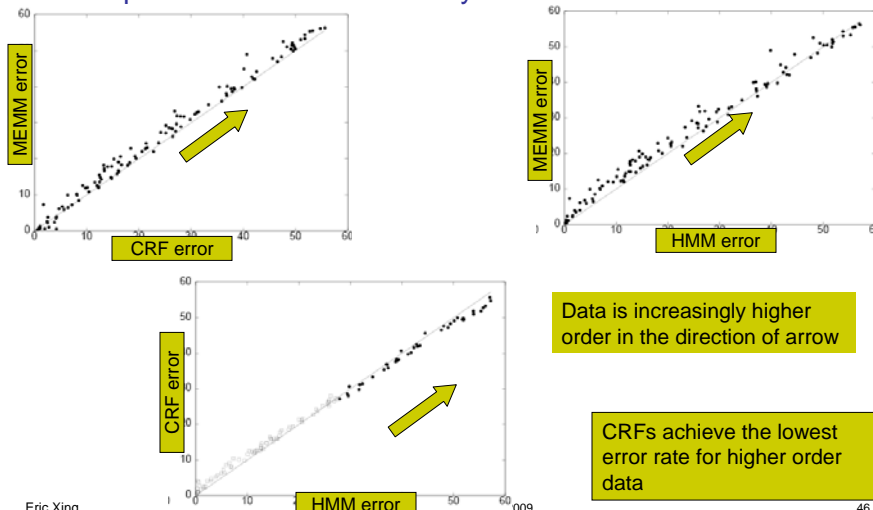
Eric Xing

© Eric Xing @ CMU, 2006-2009

45

CRFs: some empirical results

- Comparison of error rates on synthetic data



46

CRFs: some empirical results



- Parts of Speech tagging

<i>model</i>	<i>error</i>	<i>oov error</i>
HMM	5.69%	45.99%
MEMM	6.37%	54.61%
CRF	5.55%	48.05%
MEMM ⁺	4.81%	26.99%
CRF ⁺	4.27%	23.76%

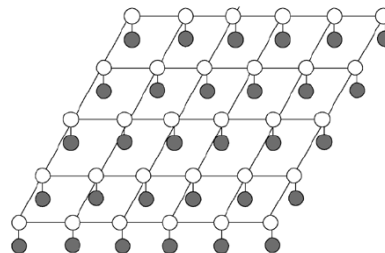
⁺Using spelling features

- Using same set of features: HMM \geq CRF > MEMM
- Using additional overlapping features: CRF⁺ > MEMM⁺ >> HMM

Other CRFs



- So far we have discussed only 1-dimensional chain CRFs
 - Inference and learning: exact
- We could also have CRFs for arbitrary graph structure
 - E.g: Grid CRFs
 - Inference and learning no longer tractable
 - Approximate techniques used
 - MCMC Sampling
 - Variational Inference
 - Loopy Belief Propagation
 - We will discuss these techniques in the future



Summary



- Conditional Random Fields are partially directed discriminative models
- Inference for 1-D chain CRFs is exact
 - Same as Max-product or Viterbi decoding
- Learning also is exact
 - globally optimum parameters can be learned
 - Requires using sum-product or forward-backward algorithm
- CRFs involving arbitrary graph structure are intractable in general
 - E.g.: Grid CRFs
 - Inference and learning require approximation techniques
 - MCMC sampling
 - Variational methods
 - Loopy BP